

PHASE EXCITED LINEAR PREDICTION ENCODER

BACKGROUND OF THE INVENTION

1. Field of the Invention

5 The present invention relates to speech coding algorithms and, more particularly to a Phase Excited Linear Predictive (PELP) low bit rate speech synthesizer and a pitch detector for a PELP synthesizer.

2. Background of Related Art

10 Mobile communications are growing at a phenomenal rate due to the success of several different second-generation digital cellular technologies, including GSM, TDMA and CDMA. To improve data throughput and sound quality, considerable effort is being devoted to the development of
15 speech coding algorithms. Indeed, speech coding is applicable to a wide range of applications, including mobile telephony, internet phones, automatic answering machines, secure speech transmission, storing and archiving speech and voice paging networks.

20 Waveform codecs are capable of providing good quality speech at bit rates down to about 16 kbits/s, but are of limited use at rates lower than 16 kbit/s. Vocoders on the other hand can provide intelligible speech at 2.4 kbits/s and below, but cannot provide natural sounding speech at
25 any bit rate. Hybrid codecs attempt to fill the gap between waveform and source codecs. The most commonly used hybrid codecs are time domain Analysis-by-Synthesis (AbS) codecs. Such codecs use the same linear prediction filter model of the vocal tract as found in Linear Predictive
30 Coding (LPC) vocoders. However, instead of applying a simple two-state, voiced/unvoiced, model to find the necessary filter input, the excitation signal is chosen by

matching the reconstructed speech waveform as closely as possible to the original speech waveform.

The distinguishing feature of AbS codecs is how the excitation waveform for the synthesis filter is chosen.

5 AbS codecs split the input speech to be coded into frames, typically about 20 ms long. For each frame, parameters are determined for a synthesis filter, and then the excitation to the synthesis filter is determined by finding the
10 filter minimizes the error between the input speech and the reconstructed speech. Thus, the encoder analyses the input speech by synthesizing many different approximations to the input speech. For each frame, the encoder transmits information representing the synthesis filter parameters
15 and the excitation to the decoder and, at the decoder, the given excitation is passed through the synthesis filter to generate the reconstructed speech. However, the numerical complexity involved in passing every possible excitation signal through the synthesis filter is quite large and
20 thus, must be reduced, but without significantly compromising the performance of the codec.

The synthesis filter is usually an all pole, short-term, linear filter intended to model the correlations introduced into speech by the action of the vocal tract.

25 The synthesis filter may also include a pitch filter to model the long-term periodicities present in voiced speech. Alternatively these long-term periodicities may be exploited by using an adaptive codebook in the excitation generator so that the excitation signal includes a
30 component of the estimated pitch period.

There are various kinds of AbS codecs, such as Multi-Pulse Excited (MPE), Regular-Pulse Excited (RPE), and Code-

Excited Linear Predictive (CELP). Generally MPE and RPE codecs will work without a pitch filter, although their performance will be improved if one is included. For CELP codecs a pitch filter is extremely important.

5 The differences between MPE, RPE and CELP codecs arise from the representation of the excitation signal. In MPE codecs, the excitation signal is given by a fixed number of non-zero pulses for every frame of speech. The positions of these non-zero pulses within the frame and their
10 amplitudes must be determined by the encoder and transmitted to the decoder. In theory it is possible to find the best values for all the pulse positions and amplitudes, but this is not practical due to the excessive complexity required. In practice some sub-optimal method
15 of finding the pulse positions and amplitudes must be used. Typically about 4 pulses per 5 ms can be used for good quality reconstructed speech at a bit-rate of around 10 kbits/s.

20 Like the MPE codec, the RPE codec uses a number of non-zero pulses to represent the excitation signal. However, the pulses are regularly spaced at a fixed interval, and the encoder only needs to determine the position of the first pulse and the amplitude of all the pulses. Therefore less information needs to be transmitted
25 about pulse positions, so for a given bit rate the RPE codec can use more non-zero pulses than the MPE codec. For example, at a bit rate of about 10 kbits/s around 10 pulses per 5 ms can be used, compared to 4 pulses for MPE codecs. This allows RPE codecs to give slightly better quality
30 reconstructed speech than MPE codecs.

Although MPE and RPE codecs provide good quality speech at rates of around 10 kbits/s and higher, they are

not suitable for lower rates due to the large amount of information that must be transmitted about the excitation pulses' positions and amplitudes. If the bit rate is reduced by using fewer pulses or by coarsely quantizing the pulse amplitudes, the reconstructed speech quality deteriorates rapidly.

Currently the most commonly used algorithm for producing good quality speech at rates below 10 kbits/s is CELP. CELP differs from MPE and RPE in that the excitation signal is effectively vector quantized. The excitation signal is given by an entry from a large vector quantizer codebook and a gain term to control its power. The codebook index is represented with about 10 bits and the gain is coded with about 5 bits. Thus, the bit rate necessary to transmit the excitation information is about 15 bits. CELP coding has been used to produce toll quality speech communications at bit rates between 4.8 and 16 kbits/s.

It is an object of the present invention to provide an efficient speech coding algorithm operable at low bit rates yet capable of reproducing high quality speech.

SUMMARY OF THE INVENTION

The present invention provides a speech encoder including a content extraction module, a pitch detector, and a naturalness enhancement module. The content extraction module includes a band pass filter that receives a speech input signal and generates a band limited speech signal. A first speech buffer connected to the band pass filter stores the band limited speech signal. An LP analysis block, connected to the first speech buffer, reads the stored speech signal and generates a plurality of LP

coefficients therefrom. An LPC to LSF block connected to the LP analysis block converts the LP coefficients to a line spectral frequency (LSF) vector. An LP analysis filter connected to the LPC to LSF block extracts an LP residual signal from the LSF vector. An LSF quantizer connected to the LPC to LSF block receives the LSF vector and determines an LSF index therefore. The pitch detector is connected to the LP analysis block of the content extraction module. The pitch detector classifies the band filtered speech signal as one of a voiced signal and an unvoiced signal. The naturalness enhancement module is connected to the content extraction module and the pitch detector. The naturalness enhancement module includes a means for extracting parameters from the LP residual signal, where for an unvoiced signal the extracted parameters include pitch and gain and for a voiced signal the extracted parameters include pitch, gain and excitation level. A quantizer quantizes the extracted parameters and generating quantized parameters.

In another embodiment, the present invention provides a content extraction module for a speech encoder. The content extraction module includes a band pass filter that receives a speech input signal and generates a band limited speech signal, and a first speech buffer connected to the band pass filter that stores the band limited speech signal. An LP analysis block connected to the first speech buffer reads the stored speech signal and generates a plurality of LP coefficients therefrom. An LPC to LSF block connected to the LP analysis block converts the LP coefficients to a line spectral frequency (LSF) vector. An LP analysis filter connected to the LPC to LSF block extracts an LP residual signal from the LSF vector, and

an LSF quantizer connected to the LPC to LSF block receives the LSF vector and determines an LSF index therefor.

In a further embodiment, the present invention provides a naturalness enhancement module for a speech encoder, where the speech encoder includes a pitch detector for determining whether an input speech signal is a voiced signal or an unvoiced signal and a content extraction module for generating an LP residual signal from the input speech signal. The naturalness enhancement module includes a means for extracting parameters from the LP residual signal, where for an unvoiced signal the extracted parameters include pitch and gain and for a voiced signal the extracted parameters include pitch, gain and excitation level, and a quantizer for quantizing the extracted parameters and generating quantized parameters.

In a further embodiment, the present invention provides a pitch detector for a speech encoder. The pitch detector includes a first operation level for analyzing a speech signal and, based on a first predetermined ambiguity value of the speech signal, generating a first estimated pitch period. A second operation level analyzes the speech signal and, based on a second predetermined ambiguity value of the speech signal, generates a second estimated pitch period.

In yet another embodiment, the present invention provides a speech signal preprocessor for preprocessing an input speech signal prior to providing the speech signal to a speech encoder. The preprocessor includes a band pass filter that receives the speech input signal and generates a band limited speech signal, and a scale down unit connected to the band pass filter for limiting a dynamic range of the band limited speech signal.

The present invention also provides a method of encoding a speech signal, including the steps of filtering the speech signal to limit its bandwidth, fragmenting the filtered speech signal into speech segments, and
5 decomposing the speech segments into a spectral envelope and an LP residual signal. The spectral envelope is represented by a plurality of LP filter coefficients (LPC). Then, the LPC are converted into a plurality of line spectral frequencies (LSF) and each speech segment is
10 classified as one of a voiced segment and an unvoiced segment based on a pitch of the segment. Next, parameters are extracted from the LP residual signal, where for an unvoiced segment the extracted parameters include pitch and gain and for a voiced segment the extracted parameters
15 include pitch, gain and excitation level. Finally, the extracted parameters are quantized to generate quantized parameters.

BRIEF DESCRIPTION OF THE DRAWINGS

20 The foregoing summary, as well as the following detailed description of preferred embodiments of the invention, will be better understood when read in conjunction with the appended drawings. For the purpose of illustrating the invention, there is shown in the drawings
25 embodiments that are presently preferred. It should be understood, however, that the invention is not limited to the precise arrangements and instrumentalities shown. In the drawings:

30 Fig. 1 is a schematic block diagram of a content extraction module of a PELP encoder in accordance with the present invention;

Fig. 2a is a schematic block diagram of a naturalness enhancement module for an unvoiced signal of a PELP encoder in accordance with the present invention;

5 Fig. 2b is a schematic block diagram of a naturalness enhancement module for a voiced signal of a PELP encoder in accordance with the present invention;

Fig. 3 is a pseudo block diagram of a pitch detector in accordance with the present invention; and

10 Fig. 4 is a flow diagram of a first PELP decoding scheme in accordance with the present invention.

DETAILED DESCRIPTION OF PREFERRED EMBODIMENTS

15 The detailed description set forth below in connection with the appended drawings is intended as a description of the presently preferred embodiments of the invention, and is not intended to represent the only forms in which the present invention may be practiced. It is to be understood that the same or equivalent functions may be accomplished by different embodiments that are intended to be
20 encompassed within the spirit and scope of the invention. In the drawings, like numerals are used to indicate like elements throughout.

The present invention is directed to a low bit rate Phase Excited Linear Predictive (PELP) speech synthesizer.
25 In PELP coding, a speech signal is classified as either voiced speech or unvoiced speech and then different coding schemes are used to process the two signals.

For voiced speech, the voiced speech signal is decomposed into a spectral envelope and a speech excitation
30 signal. An instantaneous pitch frequency is updated, for example every 5ms, to obtain a pitch contour. The pitch contour is used to extract an instantaneous pitch cycle

from the speech excitation signal. The instantaneous pitch cycle is used as a reference to extract the excitation parameters, including gain and excitation level. The spectral envelope, instantaneous pitch frequency, gains and excitation level are quantized. For unvoiced speech, a spectral envelope and gain are used, together with an unvoiced indicator.

A decoder is used to synthesize the voiced speech signal. A Linear Predictive (LP) excitation signal is constructed using a deterministic signal and a noisy signal. The LP excitation signal is then passed through a synthesis filter to generate the synthesized speech signal. To synthesize the unvoiced speech signal, a unity-power white-Gaussian noise sequence is generated and normalized to the gains to form an unvoiced excitation signal. The unvoiced excitation signal is then passed through a LP synthesis filter to generate a synthesized speech signal.

PELP coding uses linear predictive coding and mixed speech excitation to produce a natural synthesized speech signal. Different from other linear prediction based coders, the mixed speech excitation is obtained by adjusting only the phase information. The phase information is obtained using a modified speech production model. Using the modified speech production model, the information required to characterize a speech signal is reduced, which reduces the data sent over the channel. The present invention allows a natural speech signal to be synthesized with few data bits, such as at bit rates from 2.0kb/s to below 1.0kb/s.

The present invention further provides a pitch detector for the PELP coder. The pitch detector is used to classify a speech frame as either voiced or unvoiced. For

voiced speech, the pitch frequency of the voiced sound is estimated. The pitch detector is a key component of the PELP coder.

Referring now to the drawings, Figs. 1, 2a and 2b show a PELP encoder in accordance with a preferred embodiment of the present invention. The PELP encoder includes two main parts, a content extraction module **100** (Fig. 1) and a naturalness enhancement module **200a** (Fig. 2a) and **200b** (Fig. 2b).

The purpose of the content extraction module **100** is to extract the information content from an input speech signal $s'(n)$. The content extraction module **100** has a pre-processing unit that includes a band pass filter (BPF) **110**, a scale down unit **112**, and a first speech buffer **113**. The input speech signal $s'(n)$ is provided to the BPF **110**, which limits the input speech signal $s'(n)$ from about 150Hz to 3400Hz. Preferably, the BPF **110** uses an eighth order IIR filter. The aim of the lower cut-off is to reject low frequency disturbances, which could be perceptually very sensitive. The upper cut-off is to attenuate the signals at the higher frequencies. The 8th order IIR filter may be formed using a 4th order low-pass section and a 4th order high-pass section. The transfer functions of the low-pass and high-pass sections are defined in equations (1) and (2), respectively.

$$H_{lp1}(z) = \left(\frac{0.805551 + 1.611102z^{-1} + 0.805551z^{-2}}{1 + 1.518242z^{-1} + 0.703963z^{-2}} \right) \left(\frac{0.666114 + 1.332227z^{-1} + 0.666114z^{-2}}{1 + 1.255440z^{-1} + 0.409014z^{-2}} \right) \quad \text{Eqn 1}$$

$$H_{hp1}(z) = \left(\frac{0.953640 - 1.907280z^{-1} + 0.953640z^{-2}}{1 - 1.900647z^{-1} + 0.913913z^{-2}} \right) \left(\frac{0.898920 - 1.797840z^{-1} + 0.898920z^{-2}}{1 - 1.791588z^{-1} + 0.804093z^{-2}} \right) \quad \text{Eqn 2}$$

The BPF **110** thus produces a band-limited speech signal, which is provided to the scale down unit **112**. The scale

down unit **112** scales this signal down by about a half (0.5) to limit the dynamic range and hence to yield a speech signal $s(n)$. The speech signal $s(n)$ is segmented into frames, for example 20ms frames, and stored in the first speech buffer **113**. For an 8kHz sampling system, a speech frame contains 160 samples. In the presently preferred embodiment, the first speech buffer **113** stores 560 samples $B_{sp1}(n)$ for $n=0,559$ for analysis by an LP analysis block **114**. When a frame (160 samples) of the speech signal $s(n)$ is available, it is loaded into the first speech buffer **113** from samples $n=400$ to 559. The samples proceeding $B_{sp1}(400)$ are made up of the previous consecutive frames.

In the presently preferred embodiment, the LP analysis block **114** performs a 10th order Burg's LP analysis to estimate the spectral envelope of the speech frame. The LP analysis frame contains 170 samples, from $B_{sp1}(390)$ to $B_{sp1}(559)$. The result of the LP analysis is ten LP coefficients (LPC), $a''(i)$ where $i=1$ to 10. A bandwidth expansion block **116** is used to expand the set of LP coefficients using equation (3), which generates bandwidth expanded LP coefficients $a'(i)$.

$$a'(i) = 0.996^i a''(i) \quad \text{for } i=1,2,\dots,10 \quad \text{Eqn 3}$$

A frame of an LP residual signal $r(n)$ is extracted using an LP analysis filter in the following manner. After the set of bandwidth expanded LP coefficients $a'(i)$ is generated, the coefficients $a'(i)$ are converted to line spectral frequencies (LSF) $\omega'_1(i)$ ($i=1$ to 10), at an LPC to LPF block **118**. The current set of LSF $\omega'_1(i)$ is then linearly interpolated with the set of the previous frame LSF at an interpolate LSF block **120** to compute a set of

intermediate LSF $\omega_1(i)$, preferably every 5ms. Hence there are four sets of intermediate LSF $\omega_1(m,i)$ ($m=1,4$; $i=1,10$) in a speech frame. The four intermediate LSF sets $\omega_1(m,i)$ are converted back to corresponding LP coefficients $a(m,i)$ ($m=1,4$; $i=1,10$) at an LSF to LPC block **122**. Then, a frame of the residual signal $r(n)$ is obtained using an inverse filter **124** operating in accordance with equation (4).

$$r(n) = s(n) + \sum_{i=1}^{10} a(i)s(n-i) \quad \text{Eqn 4}$$

A first residual buffer **130** stores the residual signal $r(n)$. The size of the first residual buffer **130** is preferably 320 samples. That is, the stored data is $B_{rd1}(n)$ for $n=0$ to 319, which is the current residual frame and a previous consecutive frame. To compute the current residual frame, the inverse filter **124** is operated as shown in Table 1.

Filter input from $B_{sp1}(n)$ range of (n)	Filter coefficients	Filter output to $B_{rd1}(n)$ range of (n)
320 to 359	$\{a_i^{(1)}\}$	160 to 199
360 to 399	$\{a_i^{(2)}\}$	200 to 239
400 to 439	$\{a_i^{(3)}\}$	240 to 279
440 to 479	$\{a_i^{(4)}\}$	280 to 319

Table 1 Method of inverse filtering to extract excitation parameters

The LSF $\omega'_1(i)$ from the LPC to LSF block **118** are also quantized by an LSF codebook or quantizer **126** to determine an index I_L . That is, as is understood by those of ordinary skill in the art, the LSF quantizer **126** stores a number of reference LSF vectors, each of which has an index associated with it. A target LSF vector $\omega'_1(i)$ is compared with the LSF vectors stored in the LSF quantizer **126**. The

best matched LSF vector is chosen and an index I_L of the best matched LSF vector is sent over the channel for decoding.

As previously discussed, for the LP residual signal $r(n)$, different coding schemes are used for different signal types. For a voiced segment, a pitch cycle is extracted from the LP residual signal $r(n)$ every 5ms, i.e. an instantaneous pitch cycle. The gain, pitch frequency and excitation level for the instantaneous pitch cycle are extracted. A consecutive set for each parameter is arranged to form a parameter contour. The sensitivity of each parameter to the synthesised speech quality is different. Hence, different update rates are used to sample each parameter contour for coding efficiency. In the presently preferred embodiment, a 5ms update is used for gain and a 10ms update is used for the pitch frequency and excitation level. For an unvoiced segment, only the gain contour is useful. An unvoiced sub-segment is extracted from the LP residual signal $r(n)$ every 5ms. The gain of each unvoiced sub-segment is computed and arranged in time to form a gain contour. Once again a 5ms update rate is used to sample the unvoiced gain. A pitch detector **128** is used to classify the speech signal $s(n)$ as either voiced or unvoiced. In the case of voiced speech the pitch frequency is estimated.

Referring now to Fig. 3, a pseudo block diagram of the pitch detector **128** is shown. The pitch detection operation is divided into 3 levels, depending on the ambiguity of the speech signal $s(n)$.

In level (1), the speech signal $s(n)$ is filtered with a low pass filter **300** to reject the higher frequency content that may obstruct the detection of true pitch. The

cut-off frequency of the low-pass filter **300** is preferably set to 1000Hz. Preferably the filter **300** has a filter transfer function as defined in equation (5).

$$H_{lp2}(z) = \left(\frac{0.097631 + 0.195262z^{-1} + 0.097631z^{-2}}{1 - 0.942809z^{-1} + 0.333333z^{-2}} \right)$$

Eqn 5

The output $s_1(n)$ of the low-pass filter **300** is loaded into a second speech buffer **302**. In the presently preferred embodiment, the second speech buffer **302** is used to store two consecutive frames $B_{sp2}(n)$ where $n = 0$ to 319, which is 320 samples. More particularly, the input to the low pass filter **300** is taken from the first speech buffer **113** as $B_{sp1}(400)$ to $B_{sp1}(559)$ and a modified speech signal $s_1(n)$ output from the low pass filter **300** is stored in the second speech buffer **302** $B_{sp2}(160)$ to $B_{sp2}(319)$.

The stored modified speech signal $B_{sp2}(n)$, $n = 160$ to 319 is provided to an inverse filter **304** to obtain a band-limited residual signal $r_1(n)$. The filter coefficients of the inverse filter **304** are set to $a_i^{(4)}$ for $i=0,10$. The residual signal $r_1(n)$ output from the inverse filter **304** is stored in a second residual buffer **306**. The second residual buffer **306** preferably stores 320 samples $B_{rd2}(n)$ where $n=0$ to 319, and thus, the residual buffer **306** holds two consecutive residual frames. The current residual signal $r_1(n)$ is stored in $B_{rd2}(n)$, where $n=160$ to 319.

After a new residual signal $r_1(n)$ is loaded into the second residual buffer **306**, a cross-correlation function is computed at block **308** using data read from the buffer **306** $B_{rd2}(n)$ in accordance with equation (6).

$$C_r(m) = \frac{\sum_{n=319}^{160} B_{rd2}(n) B_{rd2}(n-m)}{\sqrt{\sum_{n=319}^{160} B_{rd2}^2(n) \sum_{n=319}^{160} B_{rd2}^2(n-m)}} \quad \text{Eqn 6}$$

$form=16,17,18,...,160$

A peak detector **310** finds the global maximum C_{rmax} and its location P_{rmax} , across the cross-correlation function $C_r(m)$, $m=16$ to 160 . A level detector **312** checks if C_{rmax} is greater than or equal to about 0.7 , in which case the confidence for a voice signal is high. In this case, the cross-correlation function $C_r(m)$ is re-examined to eliminate possible multiple pitch errors and hence to yield the estimated pitch-period p_{est} and its correlation function C_{est} at block **314**. The multiple-pitch error checking is preferably carried out as follows:

- i) set correlation threshold as $C_{th} = 0.75 \times C_{rmax}$
- ii) set examined range from $m = 16$ to p_{rmax}
- iii) the estimate pitch-period is equal to the first local maximum across $C_r(m)$ for $m=16$ to p_{rmax} , in ascending order of m , which has a correlation value greater than C_{th} :

$$p_{est} = Pos(C_r(p))$$

$$C_{est} = C_r(p)$$

where

$$C_r(p) \geq C_{th}$$

$$16 \leq p < P_{rmax}$$

- iv) if condition (iii) is not satisfied, then p_{est} and C_{est} are set as:

$$p_{est} = p_{rmax}$$

$$C_{est} = C_{rmax}$$

If the level detector **312** determines that C_{rmax} is less than about 0.7 , level (2) pitch detection processing is used.

Level (2)

Level (2) of the pitch detector **128** is delegated to the detection of an unvoiced signal. This is done by accessing the RMS level and energy distribution R_u of the speech signal $s(n)$. The RMS value of the speech signal $s(n)$ is computed at block **316** in accordance with equation (7).

$$RMS = \sqrt{\frac{\sum_{n=400}^{559} B_{spl}^2(n)}{160}} \quad \text{Eqn 7}$$

The vocal tract has certain major resonant frequencies that change as the configuration of the vocal tract changes, such as when different sounds are produced. The resonant peaks in the vocal tract transfer function (or frequency response) are known as "formants". It is by the formant positions that the ear is able to differentiate one speech sound from another. The energy distribution R_u , defined as the energy ratio between the higher formants and all the detectable formants, for a pre-emphasized spectral envelope, is computed at block **318**. The pre-emphasized spectral envelope is computed from a set of pre-emphasized filter coefficients that defines a system with the transfer function shown in equation (8).

$$A^{\#}(z) = (1 + 0.99z^{-1})A'(z) \quad \text{Eqn 8}$$

If a' and $a^{\#}$ are the filter coefficients for $A'(z)$ and $A^{\#}(z)$, they are related as shown in equation (9).

$$\begin{aligned} a^{\#}_0 &= 1.0 \\ a^{\#}_i &= a'_i + 0.99a'_{i-1} \quad \text{for } i = 1, 2, \dots, 10 \\ a^{\#}_{11} &= 0.99a'_{10} \end{aligned} \quad \text{Eqn 9}$$

After filter coefficients $a^{\#}$ are available, $a^{\#}$ are zero-padded to 256 samples and an FFT analysis is applied to yield a smoothed spectral envelope. For example, assuming X_k where $k=1$ to M are the magnitude values for formants (1) to (M), where formants (1) to (m) are below 2kHz and formants (m+1) to (M) are above 2kHz, the energy distribution is defined as:

$$R_u = \frac{\sum_{k=m+1}^M X_k^2}{\sum_{k=1}^M X_k^2}$$

Eqn 10

Detection of an unvoiced signal is done at block **320** by checking if either RMS is less than about 58.0 or R_u is greater than about 0.5. If either of these conditions is met, an unvoiced frame is declared and C_{est} and p_{est} are cleared or set to zero. Otherwise, the pitch detector **128** will call upon the level (3) analysis.

Level (3)

In level (3), a cross-correlation function low-pass filtered speech signal $C_s(m)$ is computed from the low-pass filtered speech signal stored in the second speech buffer **302** using equation (11), at block **322**.

$$C_s(m) = \frac{\sum_{n=319}^{160} B_{sp2}(n) B_{sp2}(n-m)}{\sqrt{\sum_{n=319}^{160} B_{sp2}^2(n) \sum_{n=319}^{160} B_{sp2}^2(n-m)}}$$

for $m = 16, 17, 18, \dots, 160$

Eqn 11

A peak detector **324** is connected to the block **322** and detects the global maximum C_{smax} and its location p_{smax} of $C_s(m)$. The correlation function $C_s(m)$ calculated at block **322** is examined at block **326**, in a similar manner as is done in level (1) with $C_r(m)$, and then the appropriate

cross-correlation function $C_r(m)$ or $C_s(m)$ is selected at block 328 to eliminate multiple pitch errors.

For example, assume the estimated pitch-period and its associated correlation function for $C_r(m)$ and $C_s(m)$ are p_{rest} and C_{rest} and p_{sest} and C_{sest} respectively. The value C_{smax} is then assessed and the following logic decisions are performed. If C_{smax} is greater than or equal to about 0.7, a voiced signal is declared and pitch logic (1) is used to choose p'_{est} from p_{rest} and p_{sest} , and determine C_{est} . The estimated pitch-period p_{est} is obtained by post processing p'_{est} . Otherwise, the sum of C_{rmax} and C_{smax} is computed, $C_{sum} = C_{rmax} + C_{smax}$. When the value of C_{sum} is available, the logic decisions are made as follows.

If $C_{sum} \geq 1.0$, a voiced signal is declared and pitch logic (2) is used to choose p'_{est} from p_{rest} and p_{sest} , and determine C_{est} . The estimated pitch-period p_{est} is obtained by post-processing p'_{est} , as described below. Otherwise, an unvoiced signal is declared, $C_{est} = 0.0$ and $p_{est} = 0$.

Pitch logic (1)

For pitch logic (1), two conditions are analyzed at a first decision block:

i) Absolute difference between the two estimated pitch periods, $p_{diff} = |p_{sest} - p_{rest}|$, is checked for $p_{diff} \geq p_{min}$, where p_{min} is a minimum pitch-period that is set to 16 samples.

ii) The value of C_{rmax} is assessed for $C_{rmax} > 0.5$.

If both conditions are met, the probability of a multiple pitch error in one of the pitch-periods (p_{sest} and p_{rest}) is high. Hence, the result is taken from the one with a smaller pitch-period:

if $p_{\text{sest}} > p_{\text{prest}}$, $p'_{\text{est}} = p_{\text{prest}}$ and $C_{\text{est}} = C_{\text{imax}}$,
 otherwise, $p'_{\text{est}} = p_{\text{sest}}$ and $C_{\text{est}} = C_{\text{smax}}$

If either of conditions (i) and (ii) fails, the results are
 5 taken from the one with a higher correlation maximum, i.e.,
 $p'_{\text{est}} = p_{\text{sest}}$ and $C_{\text{est}} = C_{\text{smax}}$.

Pitch logic (2)

Pitch logic (2) is a simple comparison between two
 10 correlation maximums. If $C_{\text{smax}} > C_{\text{imax}}$, the voicing decision
 made from $C_s(m)$ may be high, and hence the result is taken
 from $C_s(m)$, $p'_{\text{est}} = p_{\text{sest}}$ and $C_{\text{est}} = C_{\text{smax}}$. Otherwise, if $C_{\text{imax}} >$
 C_{smax} , then $p'_{\text{est}} = p_{\text{prest}}$ and $C_{\text{est}} = C_{\text{imax}}$.

After the pitch period p'_{est} is selected, the pitch
 15 period p'_{est} is smoothed by a pitch post-processing unit
330. The pitch post-processing unit **330** is a median
 smoother used to smooth out an isolated error such as a
 multiple pitch error or a sub-multiple pitch error. In the
 presently preferred embodiment, the pitch post-processing
 20 unit **330** differs from conventional median smoothers, which
 operate on the pitch-periods taken from both the previous
 and future frames, because the median smoother uses the
 current estimated pitch-period and pitch-periods estimated
 in the two previous consecutive frames.

25 Assume the estimated pitch-period for the l^{th} speech
 frame as $p(l)$ and $p(l-1)$ and $p(l-2)$ are the estimated
 pitch-periods for the two previous consecutive frames.

30 $p(l) = p'_{\text{est}}$
 $p(l-1) = p_{\text{est}}$ for $(l-1)^{\text{th}}$ frame
 $p(l-2) = p_{\text{est}}$ for $(l-2)^{\text{th}}$ frame

Three cases are analyzed.

35 i) steady voicing: $p(l) > 0$, $p(l-1) > 0$ and $p(l-2) > 0$

- ii) voice onset(2): $p(l) > 0$, $p(l-1) > 0$ and $p(l-2) = 0$
- iii) voice onset(1): $p(l) > 0$, $p(l-1) = 0$ and $p(l-2) = 0$

For steady voicing, the median smoother only operates when C_{est} is smaller than about 0.6, which is a weak voiced signal. The median smoother takes the median value of $p(l)$, $p(l-1)$ and $p(l-2)$:

$$p_{est} = \text{Median}(p(l), p(l-1), p(l-2))$$

- 10 For voice onset (2), the two estimated pitch-periods are averaged if $C_{est} < 0.5$:

$$p_{est} = 0.5 * (p(l) + p(l-1)) \quad \text{for } C_{est} < 0.5$$

- 15 This is done to ensure a smooth pitch-period trajectory. If C_{est} is greater than or equal to 0.5, a strong enough voicing can be assumed and hence $p_{est} = p(l)$. For voice onset(1), no history of pitch-periods is available and hence the estimated value is used, $p_{est} = p(l)$. Thus, the
20 pitch detector **128** indicates estimated pitch-period p_{est} and its correlation function C_{est} .

- Referring now to Figs. 2a and 2b, the naturalness enhancement module **200a/200b** of the PELP encoder is shown. In the naturalness enhancement module **200a/200b**, different
25 analyses are carried out on the residual signal $r(n)$ stored in the first residual buffer **130** (Fig. 1) for voiced and unvoiced signal types to extract a set of contours in order to enhance the quality of the synthetic speech. Fig. 2a shows the process performed on an unvoiced signal and Fig.
30 2b shows the process performed on a voiced signal.

A contour is a sequence of parameters, which in the presently preferred embodiment are updated every 5ms. As previously discussed, the length of a speech frame is 20ms, hence there are four (4) parameters (m) in a frame, which

make up a contour. The parameters for an unvoiced signal are pitch and gain. On the other hand, the parameters for a voiced signal are pitch, gain and excitation level.

5 Unvoiced signal

For an unvoiced signal, at block **210** the contours are extracted from the data $B_{rd1}(n)$ stored in the first residual buffer **130**. The contours required for an unvoiced signal are pitch and gain. The pitch contour ω_p is used to specify the pitch frequency of a speech signal at each update point. For the unvoiced signal, the pitch contour ω_p is set to zero to distinguish it from a voiced signal.

$$\omega_p(m) = 0 \quad \text{for } m=1 \text{ to } 4.$$

15 Gain factors $\lambda(m)$ are computed using the residual signal $r(n)$ data $B_{rd1}(n)$ stored in the first residual buffer **130**.

$$\lambda(m) = \sqrt{\frac{\sum_{n=n1}^{n=n1+39} B_{rd1}^2(n)}{40}}$$

Eqn 12

20 where $n1 = 160 + 40 \times (m-1)$ and $m=1$ to 4.

The encoder parameters must be quantized before being transmitted over the air to the decoder side. For the unvoiced signal, the pitch frequency and gain are quantized at block **212**, which then outputs a quantized pitch and
25 quantized gain.

Voiced signal

Three contours are required for a voiced signal, pitch, gain and excitation level. The four parameters (m)
30 for each these contours are extracted from the

instantaneous pitch cycles $u(n)$ every 5ms. Thus, at block **250** the pitch cycles $u(n)$ are extracted from the data $B_{rd1}(n)$ stored in the first residual buffer **113**. The length of each pitch cycle $u(n)$ is known as the instantaneous pitch-period $p(m)$. The value of $p(m)$ is chosen from a range of pitch-period candidates p_c . The range of p_c is computed from the estimated pitch-period p_{est} generated by the pitch detector **128**. Assume $p_c(1)$ and $p_c(M)$ are the lowest and highest pitch-period candidates, such that:

$$p_c(1) < p_c(2) < p_c(3) < \dots < p_c(M)$$

The value of $p_c(1)$ and $p_c(M)$ are computed as:

$$p_c(1) = \text{integer}(0.9 \times p_{est}) \quad \text{Eqn 13a}$$

$$p_c(M) = \text{integer}(1.1 \times p_{est}) \quad \text{Eqn 13b}$$

A cross-correlation function $C(k)$ is then computed for each of the $p_c(k)$. The $p_c(k)$ that yields the highest cross-correlation function is chosen to be the $p(m)$ at the update point. The cross-correlation function $C(k)$ is defined in equation (14).

$$C(p_{ck}) = \frac{\sum_{n=n1-1}^{n1-p_{ck}} B_{rd1}(n) B_{rd1}(n-p_{ck})}{\sqrt{\sum_{n=n1-1}^{n1-p_{ck}} B_{rd1}^2(n) \sum_{n=n1-1}^{n1-p_{ck}} B_{rd1}^2(n-p_{ck})}} \quad \text{Eqn 14}$$

The value of $n1$ is set as 200, 240, 280 and 320 for each update point. After $p(m)$ is obtained, the instantaneous pitch cycle $u(n)$ is extracted from $B_{rd1}(n)$ for the four update points.

Once an instantaneous pitch cycle $u(n)$ is available, the three contours (pitch frequency, gain and excitation level) are computed at block **252**. The gain factor λ is calculated using equation (15).

$$\lambda(m) = \sqrt{\frac{\sum_{n=0}^{p(m)-1} u(m)^2(n)}{p(m)}} \quad \text{Eqn 15}$$

To compute the excitation level ε , the absolute maximum value for the pitch cycle $u(n)$ is determined using equation (16).

$$A(m) = \max \left(|u(m, n)| \right) \quad \text{Eqn 16}$$

for $n = 0, 1, 2, \dots, p(m) - 1$

The excitation level is computed using equation (17).

$$\varepsilon(m) = 1 - \frac{\lambda(m)}{A(m)} \quad \text{Eqn 17}$$

Finally for the pitch frequency ω_p , a fractional pitch-period p' is first computed from the cross-correlation function $C(p_c(1)) \dots C(p_c(M))$. Suppose the $p(m)$ is the instantaneous pitch-period and $p(m) = p_{ck}$. The fractional pitch-period $p'(m)$ is computed as shown in equation (18).

$$p'(m) = p_{ck} + \frac{1}{2} \left(\frac{C(p_{ck} - 1) - C(p_{ck} + 1)}{C(p_{ck} - 1) - 2C(p_{ck}) + C(p_{ck} + 1)} \right) \quad \text{Eqn 18}$$

The pitch frequency is defined as shown in equation (19).

$$\omega_p(m) = \frac{2\pi}{p'(m)} \quad \text{Eqn 19}$$

Table 2 summarizes the PELP coder parameters.

Parameters	Voiced	Unvoiced
LSF	$\omega_{1i}(4) \quad i=1, 10$	$\omega_{1i}(4) \quad i=1, 10$
Gain	$\lambda(m)$	$\lambda(m)$
Pitch frequency	$\omega_p(m)$	0
Excitation level	$\varepsilon(m)$	N/A

Table 2 Summary of parameters for a PELP encoder

As with the unvoiced parameters, the encoder parameters must be quantized before being transmitted over the air to the decoder side. For the voiced signal, to achieve very low bit rate coding, at block **254**, the pitch frequency ω_p and excitation level ϵ are downsampled to reduce the information content, such as downsampling at 4:1 rate. After the pitch frequency ω_p and excitation level ϵ are downsampled, they are quantized at block **256**. Output from the quantization block **256** are a quantized pitch, quantized gain, and quantized excitation level.

Hence, only one pitch frequency and excitation level is quantized for each 20ms voiced frame. An example of the quantization scheme for a 1.8kb/s PELP coder is shown in Table 3.

Parameters	Bits/20ms frame	Method
LSF $\omega_{li}(4)$ $i=1,10$	20	Multistage-split VQ
Gain $\lambda(m)$ $m=1$ to 4	7	VQ on the logarithm gain
Pitch frequency $\omega_p(4)$	7	Scalar Quantization
Excitation level $\epsilon(4)$	2	Scalar Quantization

**Table 3 Bit allocation table for a 1.8kb/s PELP coder
(VQ – vector quantization)**

Further quality enhancement may be achieved by reducing the downsampling rate of the pitch frequency ω_p and the

- 5 excitation level ϵ , for example to 2:1 and so on, as will be understood by those of ordinary skill in the art.

PELP Decoder

The PELP decoder uses the LP residual parameters
10 generated by the encoder (gain, pitch frequency, excitation level) to reconstruct the LP excitation signal. The reconstructed LP excitation signal is a quasi-periodic signal for voiced speech and a white Gaussian noise signal for unvoiced speech. The quasi-periodic signal is
15 generated by linearly interpolating the pitch cycles at 5ms intervals. Each pitch cycle is constructed using a deterministic component and a noise component. In addition, the LSF vector is linearly interpolated with the one in the previous frame to obtain an intermediate LSF
20 vector and converted to LPC. After the excitation signal is constructed, it is passed through an LP synthesis filter to obtain the synthesised speech output signal $s(n)$.

The parameters needed for speech synthesis are listed in Table 4. If the parameters are further downsampled for
25 lower bit rates, the intermediate parameters are recovered via a linear interpolation.

PELP decoder parameters
LSF $\omega_{1,i}(4)$
Gain $\lambda(m)$
Pitch frequency $\omega_p(m)$
Excitation level $\varepsilon(m)$

Table 4 Decoder parameters

Referring now to Fig. 4, a flow diagram of a PELP decoding scheme in accordance with the present invention is shown. The speech synthesis process can be separated into two paths, one for voiced signals and one for unvoiced signals. The decision on which path to choose is based on pitch frequency ω_p . At decision block 402, if ω_p equals zero, an unvoiced signal is synthesized. On the other hand, if ω_p is greater than zero, a voiced signal is synthesized.

To synthesize an unvoiced speech frame, at block 404 a random excitation signal is generated. More particularly, four segments of a unity-power white-Gaussian sequence (40 samples each) are generated, i.e. $g'(m,n)$ for $m=1,4$; $n=0,39$. The white Gaussian noise generator is implemented by a random number generator that has a Gaussian distribution and white frequency spectrum. At block 406, each sequence $g'(m,n)$ is scaled to the corresponding gain $\lambda(m)$ to yield $g(m,n)$, as shown by equation (20).

$$\begin{aligned}
 g(m,n) &= \lambda(m)g'(m,n) \\
 &\text{for } m=1,2,3,4 \\
 &\text{for } n=0,1,2,\dots,39
 \end{aligned}$$

Eqn 20

In addition, using the codebook index I_L generated by the encode (Fig. 1) to access the LSF, four intermediate LSF vectors $\omega_{1,i}'(m,i)$ $m=1,4$; $i=1,10$ for a 20ms speech frame

are calculated at block **408**. The four intermediate LSF vectors ω_1' are then converted to LP filter coefficients $a'(m,i)$ $m=1,4$; $i=1,10$ by linearly interpolating the intermediate LSF vectors across the 20ms frame at block

- 5 **410**. More particularly, suppose the two boundary LSF vectors are $\omega_1'(l-1)$ and $\omega_1'(l)$, the LSF vector $\omega_1'(m,i)$ is then calculated as shown in equation (21).

$$\omega_1'(m,i) = \omega_1(l-1,i) + 0.25 * m * (\omega_1(l,i) - \omega_1(l-1,i)) \quad \text{Eqn 21}$$

for $i=1,2,\dots,10$

10

Finally, the synthesized unvoiced speech signal is obtained by passing the Gaussian sequence $g(m,n)$ to an LP synthesis filter **412**. The operation of the LP synthesis filter **412** is defined by difference equation (22).

15

$$s(n) = e(n) - \sum_{i=1}^{10} a_i s(n-i) \quad \text{Eqn 22}$$

where $e(n)$ is the input to the LP synthesis filter. The filtering is done according to Table 5.

Excitation signal $e(n)$	Filter coefficients	Synthesis speech $s(n)$ for $n =$
$\{g^{(1)}(n)\}$	$\{a_i^{(1)}\}$	0 to 39
$\{g^{(2)}(n)\}$	$\{a_i^{(2)}\}$	40 to 79
$\{g^{(3)}(n)\}$	$\{a_i^{(3)}\}$	80 to 119
$\{g^{(4)}(n)\}$	$\{a_i^{(4)}\}$	120 to 159

Table 5 LP synthesis filtering to generate a frame of unvoiced speech

20

A voiced speech signal is processed differently from an unvoiced speech signal. For a voiced speech signal, a quasi-periodic excitation signal is generated at block **414**. The quasi-periodic signal is generated by interpolating the

25 four synthetic pitch cycles in a 20ms frame. Each

synthetic pitch cycle is generated using the corresponding gain λ , pitch frequency ω_p and excitation level ε .

For example, suppose the synthetic pitch cycle $u(n)$ at an update point within the 20ms frame is defined in the

- 5 frequency domain by its pitch-period p , a magnitude spectrum U_k and a phase spectrum ϕ_k . Only half of the frequency spectrum is used, i.e., k is defined from $k=0$ to $k=\frac{(p+1)}{2}-1$. The pitch-period p is calculated as shown in equation (23).

10

$$p = \text{Integer} \left(\frac{2\pi}{\omega_p} \right) \quad \text{Eqn 23}$$

A flat magnitude spectrum is used in the PELP coding for U_k and is defined as shown in equation (24).

$$\begin{aligned} U_0 &= 0 \\ U_k &= \lambda \sqrt{p} \end{aligned} \quad \text{Eqn 24}$$

15

The phase spectrum ϕ_k includes deterministic phases ϕ_d at the lower frequency band and random phase components ϕ_r at the higher frequency band.

$$\phi_k = \begin{cases} \phi_{dk} & 0 < k\omega_p \leq \omega_s \\ \phi_{rk} & \omega_s < k\omega_p \leq \pi \end{cases} \quad \text{Eqn 25}$$

- 20 The separation between the two bands is known as the separation frequency ω_s , where:

$$\omega_s = \pi \times \varepsilon \quad \text{Eqn 26}$$

The deterministic phases ϕ_d are derived from a modified speech production model as shown in equation (27).

25

$$\phi_{dk} = \tan^{-1} \left(\frac{\alpha \sin(k\omega_p)}{1 - \alpha \cos(k\omega_p)} \right) + \tan^{-1} \left(\frac{\gamma \sin(k\omega_p)}{1 - \gamma \cos(k\omega_p)} \right) - 2 \tan^{-1} \left(\frac{\sin(k\omega_p)}{\beta - \cos(k\omega_p)} \right) \quad \text{Eqn 27}$$

The ways in which α , β and γ can be computed are well understood by those of ordinary skill in the art. The random phase spectrum is generated using a random number generator. The random number generator provides a uniform distributed random number range from 0 to 1.0, which is normalized to 0 and π .

After the magnitude and phase spectra for the pitch cycle are obtained, they are transformed to real and imaginary spectra for interpolation as shown in equation (28).

$$\begin{aligned} R_k &= |U_k| \cos(\phi_k) \\ I_k &= |U_k| \sin(\phi_k) \end{aligned} \quad \text{Eqn 28}$$

To synthesize a voiced excitation, the pitch frequency and the real and imaginary spectra from one pitch cycle to another are linearly interpolated to provide a smooth change of both the signal energy and shape. For example, suppose $u(m-1)(n)$ and $u(m)(n)$ are adjacent pitch cycles (5ms apart). The pitch-frequencies and real and imaginary spectra for the 2 cycles are denoted as $\omega_p(m-1)$, $R_k(m-1)$, $I_k(m-1)$ and $\omega_p(m)$, $R_k(m)$, $I_k(m)$ respectively. The voiced excitation signal $v(m)(n)$ $n=0,39$ is synthesized from these two pitch cycles using equation (29).

$$v^{(m)}(n) = \frac{1}{p^{(m)}(n)} \sum_{k=1}^{K(m)-1} \left\{ \left(R_k^{(m-1)} + \psi(n) (R_k^{(m)} - R_k^{(m-1)}) \right) \cos(k\sigma^{(m)}(n)) + \left(I_k^{(m-1)} + \psi(n) (I_k^{(m)} - I_k^{(m-1)}) \right) \sin(k\sigma^{(m)}(n)) \right\} \quad \text{Eqn 29}$$

for $n=0,1,2,\dots,39$

where $\psi(n)$ is a linear interpolation function defined by equation (30).

$$\psi(n) = \frac{n}{40} \quad \text{Eqn 30}$$

for $n = 0, 1, 2, \dots, 39$

5 The value $p(m)(n)$ is the instantaneous pitch-period for each time sample (n) , and is computed from the instantaneous pitch frequency $\omega_p(m)(n)$ as shown in equation (31).

$$p^{(m)}(n) = \frac{2\pi}{\omega_p^{(m)}(n)} \quad \text{Eqn 31}$$

10 The instantaneous pitch frequency $\omega_p^{(m)}(n)$ is computed as:

$$\omega_p^{(m)}(n) = \omega_p^{(m-1)} + \psi(n)(\omega_p^{(m)} - \omega_p^{(m-1)}) \quad \text{Eqn 32}$$

15 $K(n)$ is a parameter related to the instantaneous pitch period as:

$$K(n) = \frac{(p^{(m)}(n) + 1)}{2} \quad \text{Eqn 33}$$

The instantaneous phase value $\sigma^{(m)}(n)$ is calculated via as:

$$\sigma^{(m)}(n) = n\omega_p^{(m-1)} + \frac{n^2}{40}(\omega_p^{(m)} - \omega_p^{(m-1)}) + \sigma^{(m-1)}(40) \quad \text{Eqn 34}$$

for $n = 0, 1, 2, \dots, 39$

20 After the four pieces of voiced excitation $v(m)(n)$, $m=1, 4$; $n=0, 39$ are available, they are used as inputs to the LP synthesis filter **412** for synthesizing the voiced speech, in the same manner as is done for unvoiced speech, according
25 to Table 6.

Excitation signal $e(n)$	Filter coefficients	Synthesis speech $s(n)$ for $n =$
$\{v^{(1)}(n)\}$	$\{a'_i{}^{(1)}\}$	0 to 39
$\{v^{(2)}(n)\}$	$\{a'_i{}^{(2)}\}$	40 to 79
$\{v^{(3)}(n)\}$	$\{a'_i{}^{(3)}\}$	80 to 119
$\{v^{(4)}(n)\}$	$\{a'_i{}^{(4)}\}$	120 to 159

Table 6 LP synthesis filtering to generate a frame of voiced speech

A voiced onset frame is defined when a voiced frame is indicated directly after an unvoiced frame. In a voiced onset frame, parameters for pitch cycle $\{u^{(0)}(n)\}$ are not available for interpolating it with $\{u^{(0)}(n)\}$. To solve this problem, the parameters for $\{u^{(0)}(n)\}$ are re-used by $\{u^{(0)}(n)\}$ as shown below, and then the normal voiced synthesis is resumed.

$$p(0) = p(1)$$

$$\omega_p(0) = \omega_p(1)$$

$$R_k(0) = R_k(1)$$

$$I_k(0) = I_k(1)$$

As is apparent, the present invention provides a Phase Excited Linear Prediction type vocoder. The description of the preferred embodiments of the present invention have been presented for purposes of illustration and description, but are not intended to be exhaustive or to limit the invention to the forms disclosed. It will be appreciated by those skilled in the art that changes could be made to the embodiments described above without departing from the broad inventive concept thereof. For example, the present invention is not limited to a vocoder having any particular bit rate. It is understood, therefore, that this invention is not limited to the

particular embodiments disclosed, but covers modifications within the spirit and scope of the present invention as defined by the appended claims.

5

Table of Abbreviations and Variables

	AbS	Analysis by Synthesis
	BPF	Band Pass Filter
10	CELP	Code Excited Linear Predictive
	LP	Linear Predictive
	LPC	Linear Predictive Coefficient
	LSF	Line Spectral Frequencies
	MPE	Multi-pulse Excited
15	PELP	Phase Excited Linear Predictive
	RPE	Regular Pulse Excited
	VBR-PELP	Variable Bit Rate PELP
	$a''(i)$	LPC ($i=1,10$)
20	$a'(i)$	expanded LPC $a''(i)$
	$a(m, I)$	LPC
	$B_{sp1}(n)$	Data stored in first speech buffer 113
	$B_{sp2}(n)$	Data stored in second speech buffer 302
	$B_{rd1}(n)$	Data stored in first residual buffer 130
25	$B_{rd2}(n)$	Data stored in second residual buffer 306
	$C(k)$	cross-correlation fx for pitch period candidates
	C_{est}	cross-correlation fx of P_{est}
	$C_r(m)$	cross-correlation fx
	C_{rest}	location of P_{rest}
30	C_{rmax}	global maximum of $C_r(m)$
	$C_s(m)$	cross-correlation fx of LPF speech signal
	C_{smax}	global maximum of $C_s(m)$
	C_{sest}	location of P_{sest}
	$e(n)$	LP synthesis filter excitation signal
35	$H_{lp1}(z)$	transfer function of low pass section of BPF 110
	$H_{hp1}(z)$	transfer function of high pass section of BPF 110
	$H_{lp2}(z)$	transfer function of LPF 300
	I_L	codebook index of LSF vector $\omega_1'(i)$
	$p(m)$	instantaneous pitch period
40	p_c	pitch period candidates
	p'	fractional pitch period
	P_{est}	estimated pitch period
	P_{rest}	estimated pitch period of $C_r(m)$
	P_{rmax}	position of C_{rmax}

	P_{sest}	estimated pitch period of $C_s(m)$
	P_{smax}	position of C_{smax}
	$r(n)$	LP analysis filter residual signal
	$r_1(n)$	band limited residual signal
5	r_u	energy distribution of speech signal
	$s'(n)$	input speech signal
	$s(n)$	speech signal
	$s_1(n)$	speech signal output of LPF 300
	$u(n)$	pitch cycle
10	U_k	magnitude spectrum of pitch cycle
	$\omega_1'(i)$	LSF from $a'(i)$
	ω_1	intermediate LSF
	ω_p	pitch frequency
	λ	gain
15	ε	excitation level
	ϕ_k	phase spectrum of pitch cycle